# Writing a system call tracer using eBPF
**SH4DY** 2024-08-03 | 🏷 eBPF, linux, low-level

## # Pre-Requisites

System calls, eBPF, C, basics of low-level programming.

## # Introduction

eBPF (Extended Berkeley Packet Filter) is a technology that allows users to run custom programs within the kernel. BPF / or cBPF (classic BPF), the predecessor of eBPF provided a simple and efficient way to filter packets based on predefined rules. eBPF programs offer enhanced safety, portability, and maintainability as compared to kernel modules. There are several high-level methods available for working with eBPF programs, such as Cilium's go library, bpftrace, libbpf, etc.

- Note : This post requires the reader to have a basic understanding of eBPF . If you're not familiar with it, this post by ebpf.io is a great read.

## # Objectives

You must already be familiar with the famous tool strace . We'll be developing something similar to that using eBPF. For example,

```
1  ./beetrace /bin/ls
```

## # Concepts

Before we start writing our tool, we need to familiarize ourselves with some key concepts.

1. `Tracepoints`: They are instrumentation points placed in various parts of
   the Linux kernel code. They provide a way to hook into specific events
   or code paths within the kernel without modifying the kernel source
   code. The events available of tracing can be found at
   `/sys/kernel/debug/tracing/events`.

2. The `SEC` macro: It creates a new section with the name as the name of
   the tracepoint within the target ELF. For example,
   `SEC(tracepoint/raw_syscalls/sys_enter)` creates a new section with this
   name. The sections can be viewed using readelf.

```
1  readelf -s --wide somefile.o
```

3. `Maps`: They are shared data structures that can be accessed from both
   eBPF programs and applications running in the userspace.

# Writing the eBPF programs

We won't be writing a comprehensive tool for tracing all the system calls due
to the vast number of system calls present in the Linux kernel. Instead,
we'll focus on tracing a few common system calls. To achieve this, we'll
write two types of programs: eBPF programs and a loader (which loads the BPF
objects into the kernel and attaches them).

Let's start by creating a few data structures to set things up.

```
1   // controller.h
2
3   // SYS_ENTER : for retrieving system call arguments
4   // SYS_EXIT : for retrieving the return values of syscalls
5
6   typedef enum
7   {
8       SYS_ENTER,
9       SYS_EXIT
10  } event_mode;
11
12  struct inner_syscall_info
13  {
14      union
```

```c
15        {
16            struct
17            {
18                // For SYS_ENTER mode
19                char name[32];
20                int num_args;
21                long syscall_nr;
22                void *args[MAX_ARGS];
23            };
24            long retval; // For SYS_EXIT mode
25        };
26        event_mode mode;
27    };
28
29    struct default_syscall_info{
30        char name[32];
31        int num_args;
32    };
33
34    // Array for storing the name and argument count of system calls
35    const struct default_syscall_info syscalls[MAX_SYSCALL_NR] = {
36        [SYS_fork] = {"fork", 0},
37        [SYS_alarm] = {"alarm", 1},
38        [SYS_brk] = {"brk", 1},
39        [SYS_close] = {"close", 1},
40        [SYS_exit] = {"exit", 1},
41        [SYS_exit_group] = {"exit_group", 1},
42        [SYS_set_tid_address] = {"set_tid_address", 1},
43        [SYS_set_robust_list] = {"set_robust_list", 1},
44        [SYS_access] = {"access", 2},
45        [SYS_arch_prctl] = {"arch_prctl", 2},
46        [SYS_kill] = {"kill", 2},
47        [SYS_listen] = {"listen", 2},
48        [SYS_munmap] = {"sys_munmap", 2},
49        [SYS_open] = {"open", 2},
50        [SYS_stat] = {"stat", 2},
51        [SYS_fstat] = {"fstat", 2},
52        [SYS_lstat] = {"lstat", 2},
53        [SYS_accept] = {"accept", 3},
54        [SYS_connect] = {"connect", 3},
55        [SYS_execve] = {"execve", 3},
56        [SYS_ioctl] = {"ioctl", 3},
57        [SYS_getrandom] = {"getrandom", 3},
58        [SYS_lseek] = {"lseek", 3},
59        [SYS_poll] = {"poll", 3},
60        [SYS_read] = {"read", 3},
61        [SYS_write] = {"write", 3},
```

```
62        [SYS_mprotect] = {"mprotect", 3},
63        [SYS_openat] = {"openat", 3},
64        [SYS_socket] = {"socket", 3},
65        [SYS_newfstatat] = {"newfstatat", 4},
66        [SYS_pread64] = {"pread64", 4},
67        [SYS_prlimit64] = {"prlimit64", 4},
68        [SYS_rseq] = {"rseq", 4},
69        [SYS_sendfile] = {"sendfile", 4},
70        [SYS_socketpair] = {"socketpair", 4},
71        [SYS_mmap] = {"mmap", 6},
72        [SYS_recvfrom] = {"recvfrom", 6},
73        [SYS_sendto] = {"sendto", 6},
74    };
```

The loader will read the path of the ELF file to be traced, which will be provided by the user as a command line argument. Then, the loader will spawn a child process and use `execve` to run the program specified in the command line argument.

The parent process will handle all the necessary setup for loading and attaching the eBPF programs. It also performs the crucial task of sending the child process's ID to the eBPF program via the BPF hashmap.

```c
1   // loader.c
2
3   int main(int argc, char **argv)
4   {
5     if (argc < 2)
6     {
7       fatal_error("Usage: ./beetrace <path_to_program>");
8     }
9
10    const char *file_path = argv[1];
11
12    pid_t pid = fork();
13    if (pid == 0)
14    {
15      // Child process
16      int fd = open("/dev/null", O_WRONLY);
17      if(fd==-1){
18          // error
19      }
20      dup2(fd, 1); // disable stdout for the child process
21      sleep(2); // wait for the parent process to do the required setup for t
22      execve(file_path, NULL, NULL);
```

```
23      }
24    else{
25        // Parent process
26      }
27  }
```

To trace system calls, we need to write eBPF programs that are triggered by the `tracepoint/raw_syscalls/sys_enter` and `tracepoint/raw_syscalls/sys_exit` tracepoints. These tracepoints provide access to the system call number and arguments. For a given system call, the `tracepoint/raw_syscalls/sys_enter` tracepoint is always triggered before the `tracepoint/raw_syscalls/sys_exit` tracepoint. We can use the former to retrieve the system call arguments and the latter to obtain the return value. Additionally, we will use eBPF maps to share information between the user-space program and our eBPF programs. Specifically, we will use two types of eBPF maps: hashmaps and ring buffers.

```c
1   // controller.c
2
3   // Hashmap
4   struct
5   {
6     __uint(type, BPF_MAP_TYPE_HASH);
7     __uint(key_size, 10);
8     __uint(value_size, 4);
9     __uint(max_entries, 256 * 1024);
10  } pid_hashmap SEC(".maps");
11
12  // Ring buffer
13  struct
14  {
15    __uint(type, BPF_MAP_TYPE_RINGBUF);
16    __uint(max_entries, 256 * 1024);
17  } syscall_info_buffer SEC(".maps");
18
```

Having defined the maps, we're ready to write the programs. Let's start by writing the program for the tracepoint `tracepoint/raw_syscalls/sys_enter`.

```c
1   // loader.c
2
3   SEC("tracepoint/raw_syscalls/sys_enter")
4   int detect_syscall_enter(struct trace_event_raw_sys_enter *ctx)
```

```c
 5  {
 6      // Retrieve the system call number
 7      long syscall_nr = ctx->id;
 8      const char *key = "child_pid";
 9      int target_pid;
10
11      // Reading the process id of the child process in userland
12      void *value = bpf_map_lookup_elem(&pid_hashmap, key);
13      void *args[MAX_ARGS];
14
15      if (value)
16      {
17          target_pid = *(int *)value;
18
19          // PID of the process that executed the current system call
20          pid_t pid = bpf_get_current_pid_tgid() & 0xffffffff;
21          if (pid == target_pid && syscall_nr >= 0 && syscall_nr < MAX_SYSCALL_NR
22          {
23
24              int idx = syscall_nr;
25              // Reserve space in the ring buffer
26              struct inner_syscall_info *info = bpf_ringbuf_reserve(&syscall_info_b
27              if (!info)
28              {
29                  bpf_printk("bpf_ringbuf_reserve failed");
30                  return 1;
31              }
32
33              // Copy the syscall name into info->name
34              bpf_probe_read_kernel_str(info->name, sizeof(syscalls[syscall_nr].nam
35              for (int i = 0; i < MAX_ARGS; i++)
36              {
37                  info->args[i] = (void *)BPF_CORE_READ(ctx, args[i]);
38              }
39              info->num_args = syscalls[syscall_nr].num_args;
40              info->syscall_nr = syscall_nr;
41              info->mode = SYS_ENTER;
42              // Insert into ring buffer
43              bpf_ringbuf_submit(info, 0);
44          }
45      }
46      return 0;
47  }
```

Similarly, we can write the program for reading the return value and sending it to userland.

```c
// controller.c

SEC("tracepoint/raw_syscalls/sys_exit")
int detect_syscall_exit(struct trace_event_raw_sys_exit *ctx)
{
  const char *key = "child_pid";
  void *value = bpf_map_lookup_elem(&pid_hashmap, key);
  pid_t pid, target_pid;

  if (value)
  {
    pid = bpf_get_current_pid_tgid() & 0xffffffff;
    target_pid = *(pid_t *)value;
    if (pid == target_pid)
    {
      struct inner_syscall_info *info = bpf_ringbuf_reserve(&syscall_info_b
      if (!info)
      {
        bpf_printk("bpf_ringbuf_reserve failed");
        return 1;
      }
      info->mode = SYS_EXIT;
      info->retval = ctx->ret;
      bpf_ringbuf_submit(info, 0);
    }
  }
  return 0;
}
```

Let's now finalize the functionality for the parent process in the loader program. Before doing that, we need to understand how some key functions work.

1. `bpf_object__open`: Creates a bpf_object by opening the BPF ELF object file pointed to by the passed path and loading it into memory.

```c
LIBBPF_API struct bpf_object *bpf_object__open(const char *path);
```

2. `bpf_object__load`: Loads BPF object into kernel.

```
1   LIBBPF_API int bpf_object__load(struct bpf_object *obj);
```

3. `bpf_object__find_program_by_name` : Returns a pointer to a valid BPF program.

```
1   LIBBPF_API struct bpf_program *bpf_object__find_program_by_name(const struct
```

4. `bpf_program__attach` : Function for attaching a BPF program based on auto-detection of program type, attach type, and extra paremeters, where applicable.

```
1   LIBBPF_API struct bpf_link *bpf_program__attach(const struct bpf_progra
```

5. `bpf_map__update_elem` : Allows to insert or update value in BPF map that corresponds to provided key.

```
1   LIBBPF_API int bpf_map__update_elem(const struct bpf_map *map,const voi
```

6. `bpf_object__find_map_fd_by_name` : Given a BPF map name, it returns a file descriptor to it.

```
1   LIBBPF_API int bpf_object__find_map_fd_by_name(const struct bpf_object
```

7. `ring_buffer__new` : Returns a pointer to the ring buffer.

```
1   LIBBPF_API struct ring_buffer *ring_buffer__new(int map_fd, ring_buffer
```

The second argument must be a function which can be used for handling the data received from the ring buffer.

```
1   bool initialized = false;
2
3   static int syscall_logger(void *ctx, void *data, size_t len)
4   {
```

```
 5      struct inner_syscall_info *info = (struct inner_syscall_info *)data;
 6      if (!info)
 7      {
 8        return -1;
 9      }
10
11      if (info->mode == SYS_ENTER)
12      {
13        initialized = true;
14        printf("%s(", info->name);
15        for (int i = 0; i < info->num_args; i++)
16        {
17          printf("%p,", info->args[i]);
18        }
19        printf("\b) = ");
20      }
21      else if (info->mode == SYS_EXIT)
22      {
23        if (initialized)
24        {
25          printf("0x%lx\n", info->retval);
26        }
27      }
28      return 0;
29    }
30
```

It prints the name and arguments of the system calls.

8. `ring_buffer__consume`: It processes the available events in the ring
   buffer.

```
1   LIBBPF_API int ring_buffer__consume(struct ring_buffer *rb);
```

We now have everything needed to write the loader.

```
1   // loader.c
2   #include <bpf/libbpf.h>
3   #include "controller.h"
4   #include <fcntl.h>
5   #include <stdio.h>
6   #include <stdlib.h>
7   #include <sys/wait.h>
```

```c
 8   #include <unistd.h>

 9

10   void fatal_error(const char *message)
11   {
12     puts(message);
13     exit(1);
14   }

15

16   bool initialized = false;

17

18   static int syscall_logger(void *ctx, void *data, size_t len)
19   {
20     struct inner_syscall_info *info = (struct inner_syscall_info *)data;
21     if (!info)
22     {
23       return -1;
24     }

25

26     if (info->mode == SYS_ENTER)
27     {
28       initialized = true;
29       printf("%s(", info->name);
30       for (int i = 0; i < info->num_args; i++)
31       {
32         printf("%p,", info->args[i]);
33       }
34       printf("\b) = ");
35     }
36     else if (info->mode == SYS_EXIT)
37     {
38       if (initialized)
39       {
40         printf("0x%lx\n", info->retval);
41       }
42     }
43     return 0;
44   }

45

46   int main(int argc, char **argv)
47   {
48     int status;
49     struct bpf_object *obj;
50     struct bpf_program *enter_prog, *exit_prog;
51     struct bpf_map *syscall_map;
52     const char *obj_name = "controller.o";
53     const char *map_name = "pid_hashmap";
54     const char *enter_prog_name = "detect_syscall_enter";
```

```c
55      const char *exit_prog_name = "detect_syscall_exit";
56      const char *syscall_info_bufname = "syscall_info_buffer";
57
58      if (argc < 2)
59      {
60        fatal_error("Usage: ./beetrace <path_to_program>");
61      }
62      const char *file_path = argv[1];
63
64      pid_t pid = fork();
65      if (pid == 0)
66      {
67        int fd = open("/dev/null", O_WRONLY);
68        if(fd==-1){
69          fatal_error("failed to open /dev/null");
70        }
71        dup2(fd, 1);
72        sleep(2);
73        execve(file_path, NULL, NULL);
74      }
75      else
76      {
77        printf("Spawned child process with a PID of %d\n", pid);
78        obj = bpf_object__open(obj_name);
79        if (!obj)
80        {
81          fatal_error("failed to open the BPF object");
82        }
83        if (bpf_object__load(obj))
84        {
85          fatal_error("failed to load the BPF object into kernel");
86        }
87
88        enter_prog = bpf_object__find_program_by_name(obj, enter_prog_name);
89        exit_prog = bpf_object__find_program_by_name(obj, exit_prog_name);
90
91        if (!enter_prog || !exit_prog)
92        {
93          fatal_error("failed to find the BPF program");
94        }
95        if (!bpf_program__attach(enter_prog) || !bpf_program__attach(exit_prog
96        {
97          fatal_error("failed to attach the BPF program");
98        }
99        syscall_map = bpf_object__find_map_by_name(obj, map_name);
100       if (!syscall_map)
101       {
```

```
102        fatal_error("failed to find the BPF map");
103      }
104      const char *key = "child_pid";
105      int err = bpf_map__update_elem(syscall_map, key, 10, (void *)&pid, siz
106      if (err)
107      {
108        printf("%d", err);
109        fatal_error("failed to insert child pid into the ring buffer");
110      }
111
112      int rbFd = bpf_object__find_map_fd_by_name(obj, syscall_info_bufname);
113
114      struct ring_buffer *rbuffer = ring_buffer__new(rbFd, syscall_logger, N
115
116      if (!rbuffer)
117      {
118        fatal_error("failed to allocate ring buffer");
119      }
120
121      if (wait(&status) == -1)
122      {
123        fatal_error("failed to wait for the child process");
124      }
125
126      while (1)
127      {
128        int e = ring_buffer__consume(rbuffer);
129        if (!e)
130        {
131          break;
132        }
133        sleep(1);
134      }
135    }
136    return 0;
137  }
```

And, here are the eBPF programs. The C code will be compiled into a single
object file.

```
1  // controller.c
2
3  #include "vmlinux.h"
4  #include <bpf/bpf_helpers.h>
```

```
5   #include <bpf/bpf_core_read.h>
6   #include <sys/syscall.h>
7   #include "controller.h"
8
9   struct
10  {
11    __uint(type, BPF_MAP_TYPE_HASH);
12    __uint(key_size, 10);
13    __uint(value_size, 4);
14    __uint(max_entries, 256 * 1024);
15  } pid_hashmap SEC(".maps");
16
17  struct
18  {
19    __uint(type, BPF_MAP_TYPE_RINGBUF);
20    __uint(max_entries, 256 * 1024);
21  } syscall_info_buffer SEC(".maps");
22
23
24  SEC("tracepoint/raw_syscalls/sys_enter")
25  int detect_syscall_enter(struct trace_event_raw_sys_enter *ctx)
26  {
27    // Retrieve the system call number
28    long syscall_nr = ctx->id;
29    const char *key = "child_pid";
30    int target_pid;
31
32    // Reading the process id of the child process in userland
33    void *value = bpf_map_lookup_elem(&pid_hashmap, key);
34    void *args[MAX_ARGS];
35
36    if (value)
37    {
38      target_pid = *(int *)value;
39
40      // PID of the process that executed the current system call
41      pid_t pid = bpf_get_current_pid_tgid() & 0xffffffff;
42      if (pid == target_pid && syscall_nr >= 0 && syscall_nr < MAX_SYSCALL_NR
43      {
44
45        int idx = syscall_nr;
46        // Reserve space in the ring buffer
47        struct inner_syscall_info *info = bpf_ringbuf_reserve(&syscall_info_b
48        if (!info)
49        {
50          bpf_printk("bpf_ringbuf_reserve failed");
51          return 1;
```

```
52            }
53
54            // Copy the syscall name into info->name
55            bpf_probe_read_kernel_str(info->name, sizeof(syscalls[syscall_nr].nam
56            for (int i = 0; i < MAX_ARGS; i++)
57            {
58              info->args[i] = (void *)BPF_CORE_READ(ctx, args[i]);
59            }
60            info->num_args = syscalls[syscall_nr].num_args;
61            info->syscall_nr = syscall_nr;
62            info->mode = SYS_ENTER;
63            // Insert into ring buffer
64            bpf_ringbuf_submit(info, 0);
65          }
66        }
67      return 0;
68  }
69
70  SEC("tracepoint/raw_syscalls/sys_exit")
71  int detect_syscall_exit(struct trace_event_raw_sys_exit *ctx)
72  {
73      const char *key = "child_pid";
74      void *value = bpf_map_lookup_elem(&pid_hashmap, key);
75      pid_t pid, target_pid;
76
77      if (value)
78      {
79        pid = bpf_get_current_pid_tgid() & 0xffffffff;
80        target_pid = *(pid_t *)value;
81        if (pid == target_pid)
82        {
83          struct inner_syscall_info *info = bpf_ringbuf_reserve(&syscall_info_b
84          if (!info)
85          {
86            bpf_printk("bpf_ringbuf_reserve failed");
87            return 1;
88          }
89          info->mode = SYS_EXIT;
90          info->retval = ctx->ret;
91          bpf_ringbuf_submit(info, 0);
92        }
93      }
94      return 0;
95  }
96
97  char LICENSE[] SEC("license") = "GPL";
98
```

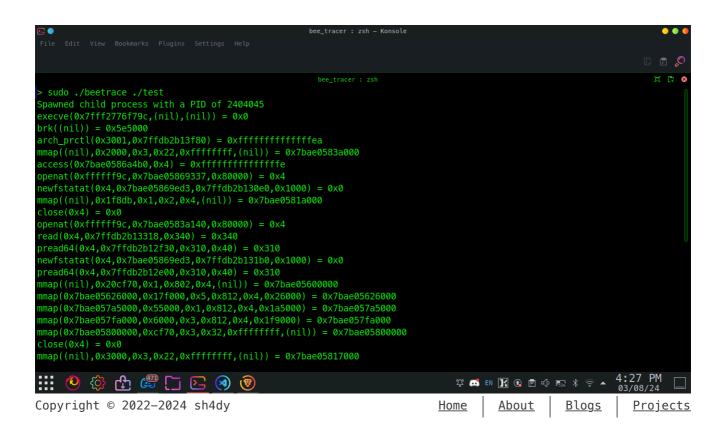Before compiling, we can create a test program which will be traced by our tool.

```
1   #include<stdio.h>
2   int main(){
3       puts("tracer in action");
4       return 0;
5   }
```

The following Makefile can be used to compile all the stuff.

```
1   compile:
2           clang -O2 -g -Wall -I/usr/include -I/usr/include/bpf -o beetrace loa
3           clang -O2 -g -target bpf -c controller.c -o controller.o
4
```

Now let's execute the loader with root privileges.

```
1   sudo ./beetrace ./test
```

```
bee_tracer : zsh — Konsole
File  Edit  View  Bookmarks  Plugins  Settings  Help

                              bee_tracer : zsh

> sudo ./beetrace ./test
Spawned child process with a PID of 2404045
execve(0x7fff2776f79c,(nil),(nil)) = 0x0
brk((nil)) = 0x5e5000
arch_prctl(0x3001,0x7ffdb2b13f80) = 0xffffffffffffffea
mmap((nil),0x2000,0x3,0x22,0xffffffff,(nil)) = 0x7bae0583a000
access(0x7bae0586a4b0,0x4) = 0xfffffffffffffffe
openat(0xffffff9c,0x7bae05869337,0x80000) = 0x4
newfstatat(0x4,0x7bae05869ed3,0x7ffdb2b130e0,0x1000) = 0x0
mmap((nil),0x1f8db,0x1,0x2,0x4,(nil)) = 0x7bae0581a000
close(0x4) = 0x0
openat(0xffffff9c,0x7bae0583a140,0x80000) = 0x4
read(0x4,0x7ffdb2b13318,0x340) = 0x340
pread64(0x4,0x7ffdb2b12f30,0x310,0x40) = 0x310
newfstatat(0x4,0x7bae05869ed3,0x7ffdb2b131b0,0x1000) = 0x0
pread64(0x4,0x7ffdb2b12e00,0x310,0x40) = 0x310
mmap((nil),0x20cf70,0x1,0x802,0x4,(nil)) = 0x7bae05600000
mmap(0x7bae05626000,0x17f000,0x5,0x812,0x4,0x26000) = 0x7bae05626000
mmap(0x7bae057a5000,0x55000,0x1,0x812,0x4,0x1a5000) = 0x7bae057a5000
mmap(0x7bae057fa000,0x6000,0x3,0x812,0x4,0x1f9000) = 0x7bae057fa000
mmap(0x7bae05800000,0xcf70,0x3,0x32,0xffffffff,(nil)) = 0x7bae05800000
close(0x4) = 0x0
mmap((nil),0x3000,0x3,0x22,0xffffffff,(nil)) = 0x7bae05817000
```

Copyright © 2022-2024 sh4dy                           Home  |  About  |  Blogs  |  Projects

The entire code can be found in this GitHub repository.


References:


https://ebpf.io/


https://github.com/libbpf/libbpf

The entire code can be found in this GitHub repository.


References: